

# Semantic Formalization of Cross-site User Browsing Behavior

Julia Hoxha  
 Karlsruhe Institute of Technology  
 Institute of Applied Informatics and Formal Description Methods (AIFB)  
 Karlsruhe, Germany  
 julia.hoxha@kit.edu

## Keywords

semantic log, cross-site browsing log formalization, classification of navigation logs, prediction with structured output, structured SVM

## 1. INTRODUCTION

Large amounts of data are being produced daily as detailed records of Web usage behavior, but the task of deriving knowledge from them still remains a challenge. Modeling and mining approaches are significant instruments to discover browsing patterns in such data and to understand how users browse Web sites.

There is an increasing body of literature on the investigation of clickstream data and navigation behavior modeling, with the majority focusing on data collected in a single site. Inspiring works [25] convincingly argue on the benefits of studying user behavior at multiple websites. Such approaches present significant potential to derive actionable behavioral knowledge and make better future forecasts, but they still have to tackle the problem of heterogeneity of the information encountered at different sites.

We approach the problem of usage data comprehensibility at its root, addressing the issue of semantically formalizing cross-site user Web browsing behavior. Usage data (or usage logs) are syntactic representations of Uniform Resource Locator (URL) requests of pages and Web resources accessed by the site visitors. Due to the primarily syntactical nature of such requests, comprehension of users' browsing patterns is difficult. Hence, there is an urge for formalization approaches that leverage the semantics of the usage data in accordance with the domain they occurred.

As such, mapping usage logs to comprehensible events from the application domain helps to discover more insights about user behavior. While most approaches use flat taxonomies to represent such vocabulary, we deploy ontologies for structuring domain concepts and relations, since they ensure a

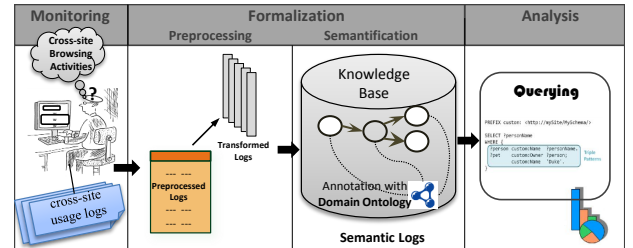


Figure 1: User Browsing Behavior Formalization Approach

richer semantic model of a Web site content.

This work aims at monitoring user behavior across multiple Web sites, logging clickthrough data upon agreement of Internet users. Each log entry is a tuple  $\mathcal{L} = \{UserID, URL, timestamp\}$  of a user ID, URL of the accessed Web resource, and real time when this happened. These usage data logs are initially stored in raw form, as produced upon each user interaction.

The overall approach, illustrated in Fig. 1, comprises a series of steps, such as data preprocessing (human logs filtering, session construction, data transformation), formalization of usage logs, and techniques for their semantic enrichment.

This work gives the following contributions to the field:

- **Model for the formal and semantic representation of cross-site browsing logs.** I present the Web browsing Activity Model (WAM), expressed as an OWL-2-DL ontology, which enables a shared conceptualization of the knowledge from the various domains where the usage logs are recorded.
- **Techniques for the automatic extraction of the usage logs semantics** from heterogeneous sources, in which the domain knowledge has a semi-structured or structured formal representation.
- **New approach for semi-supervised prediction with ontology-based output spaces.** This covers the problem of inferring the semantics of logs belonging to sites that do not offer a domain ontology. The contribution is a structured prediction algorithm formulated for the case of complex output objects rep-

resented as ontologies (with is-a hierarchies of classes and relations among them).

## 2. MOTIVATION

Existing approaches can benefit from leveraging usage data with semantics in the following ways: increase understandability of user behavior with respect to the application domain; enable analysis on higher levels of abstraction e.g. for parameters in URL (museum instead of Louvre), or for location (capital instead of Paris), which can be also used for privacy protection; allow formulation of more expressive queries for mining user behavioral patterns. Furthermore, enrichment of usage data with domain knowledge provides a broader context of user behavior, which can be exploited for more intelligent recommendation models.

The semantically-leveraged logs provide an added-value with respect to their syntactic representation in being useful inputs for techniques such as semantic pattern mining, next-step navigation prediction or user clustering, which usually assume that the semantics of logs exists or are manually derived. A more beneficial aspect is the extension of these techniques to deal with cross-site browsing data and not only a single Web site.

### 2.1 Applications

An interesting application is the integration of domain knowledge in the process of discovering usage patterns. This helps to increase the precision, and hence interpretation of the retrieved patterns, while ensuring different level of abstraction. I present two approaches for discovering browsing behavior patterns, while using as basis the formal and semantic representation of logs:

#### I. Ontology-based Web usage mining

The first application deals with the automatic mining of frequent patterns from the sequence of event logs, which are enriched with description from domain ontologies. While recent trends in Web Usage Mining (WUM) have put the emphasis on the exploitation of ontologies to the pattern mining process, yet they share two limitations: the ontologies are either restricted to representations of class taxonomies while ignoring relates among the concepts, which reduces the problem back to the traditional generalized sequential pattern discovery, or they are restricted to a single ontology (single Site) that is assumed to be completed with relations. Because of the heterogeneity of Web sites and respective domain knowledge, our setting requires a mining technique that addresses the problem when there are multiple ontologies in background and not all the relations among the concepts are established. Hence, they still need to be inferred during the mining process.

As a contribution to the WUM field, with practical motivation from the Web personalization field, I propose an approach for mining frequent sequential patterns in the presence of multiple domain ontologies. The mined patterns can, then, easily be extended to association rules, which provides predictions for the user's next step navigation preference.

#### II. Pattern discovery with $\mathcal{DL}$ -LTL expressive queries

In this application, patterns are discovered from the corpus of the semantically formalized logs upon issuing specific

queries that express semantic and temporal conditions of usage behavior.

While the first application (mining) concentrates on the semantics of the logs, another crucial aspect to consider when analyzing browsing behavior is also its temporal dynamic. Additional aspects of user browsing behavior can be discovered if reasoning not only with semantic constraints, but also with more expressive temporal conditions is made possible. I introduce an approach to formulate queries using a temporalized description logic called  $\mathcal{DL}$ -LTL, which combines *SROIQ* [12] with Lineal Temporal Logic (LTL) [2] over finite traces.

It is further shown how to search for behavioral patterns from the usage logs applying a query answering technique, which is based on current model checking tools. This allows to automatically retrieve sessions of user browsing events that satisfy a set of semantic and temporal conditions. The adaptation and application of the  $\mathcal{DL}$ -LTL logic and these techniques for the setting of Web usage analysis are novel.

## 3. STATE OF THE ART

The contributions related to this work are grouped into works dealing with 1) the modeling of user browsing behavior at multiple Web sites, 2) formal and semantic description of usage logs, and 3) exploitation of ontologies in Web usage mining, and 4) prediction of structured data.

### 3.1 Modeling Cross-site User Browsing Behavior

Interest to characterize online behavior has started much earlier with works such as those of Catledge *et al.* [6], and Montgomery *et al.* [23] that try to identify browsing strategies and patterns in the web. Browsing activity has been studied and modeled, e.g. Bucklin *et al.* [5] and others, usually exploiting server-side logs of visitors in a specific website.

Regarding the modeling of browsing behavior at multiple websites, Downey *et al.* [10] propose a state machine representation for describing search activities. They present an approach for modeling and analyzing user behavior, focusing on the search activities and what users do when they depart the search engine. Park and Fader [25] present a stochastic timing model of cross-site user visit behavior, using information from one site to explain the behavior at another. While, Johnson *et al.* [16] study online search and browsing behavior across competing e-commerce sites.

The works in this category do not particularly apply semantic techniques or ontologies for behavior modeling.

### 3.2 Ontologies in Usage Mining

There is an extensive body of work dealing with usage log analysis and mining, but we focus on the combination of these techniques with semantic technologies, which start with contributions such as Stumme *et al.* [28] and Oberle *et al.* [24]. In this field, research has been mostly focused on search query logs or user profiling. Recent approaches, which use semantics for extracting behavior patterns from

web navigation logs, are presented by Yilmaz *et al.* [34] and Mabroukeh *et al.* [19].

Vanzin *et al.* [32] present ontology-based filtering mechanisms for the retrieval of Web usage patterns. More recently, Mehdi *et al.* [21] tackle the problem of mining meaningful usage patterns and exploit the impact of ontologies to solve this problem. These works are restricted to only one domain and not cross-site browsing behavior. Hence, they mostly deal with a mining problem in the presence of a single ontology. It is interesting to explore further the discovery of patterns when multiple domain ontologies are involved, considering the establishment of mappings between them as an additional requirement of the mining process.

It is important to note though, that the process of enriching of logs with semantics is not the central problem of these works. They mostly use the ontological knowledge in the background for leveraging or optimizing the mining techniques.

### 3.3 Semantic Formalization of Usage Logs

This group consists of works that directly deal with semantic annotation of usage logs, hence mapping the requests of Web resources to meaningful concepts from the application domain. d'Aquin *et al.* [9] present the UCIAD platform<sup>1</sup>, which applies annotation of user-centric activity data. It relies on pre-defined URL patterns to characterize accessed resources over which the activities are realised, and therefore their respective semantics. As part of setting up the platform, it is initially defined which is the set of websites that are present on the considered server, as well as the URL patterns, expressed as regular expressions, enable to recognise webpages as parts of these websites. Similarly, definitions of the user activities are also manually made in the setup process, in order to characterize and give semantics to the user actions.

The work of Tvarozek *et al.* [31], while actually focusing on an architecture for the personalized presentation layer of Web-based information systems, covers in one of its techniques the problem of semantically annotating usage logs. In order to create comprehensive logs of user actions, the logs browsing events captured by a client side monitoring tool, as well as server-side logging data, are enhanced with semantics from the Web sites content using a SemanticLog tool. This tool is based on a semantically-enabled portal, which means that there is a conceptual ontology in the background of the site. The mapping of an interaction of the user with parts of the Web site, then use this ontology to generate the annotation of the user action. In this case, the semantics of the logs are not inferred, but rather defined in background as part of the engineering of the site. Still, this can be feasible only in the case when one is in charge of the content of the site, and also restricted to a small set of sites. Additional manual effort in the engineering process is needed to generate the semantic annotations.

Stühmer *et al.* [27] focus on processing complex events of user interactions with annotated Web pages, and they also present an approach for capturing and lifting these events

<sup>1</sup><http://uciad.info/ub/>

in RDF. Hence, instead of dealing with the syntactical form of events, they also address leveraging logs with semantic information, which pertains to the actual domain knowledge of the Web page. As in the previous work, this technique also assumes the presence of a semantically-enabled Web site. In this case, RDFa is used to support the semantics embedded within actual Web page data and allow reusable semantic markup inside of Web pages.

The related works in this group are restricted to a manual approach for enriching the logs with semantics. This limitation poses a significant burden when we need to analyse browsing behavior at various Web sites, which leads to immense efforts of extracting the semantics of logs and mapping them to respective domain ontologies. Moreover, it is assumed that the domain ontology is provided. This leaves the problem of inferring (learning) the semantic types of logs for non-semantically enabled sites still a challenge.

### 3.4 Prediction of Structured Data

Machine Learning today offers a broad range of methods for classification and regression, but only a few cover the problem of predicting complex objects, such as trees or graphs. The approaches dealing with prediction of structured and interdependent output data are principally grouped into those using probabilistic models (e.g. Conditional Graphical Models, HMM) and those using discriminative models (e.g. Max-Margin Structured Classification, Energy-Based Models, SVM).

In the latter group, Support Vector Machines (SVM) for structured and interdependent output spaces [30, 15] offer solid theoretical foundations, as well as very high efficiency for the structured prediction approach. While structural SVMs provide a generalized formulation of the learning problem, its state of the art applications cover only the case when the output object are sequences or trees.

There is still the need to reformulate the learning problem, and further adapt the SVMs for the case when the output instances are objects represented as ontologies. In this case, ontologies comprise not only a hierarchical structure of the classes (is-a hierarchy) in the output space, but also a set of semantic relations between these classes. The difficulty of the prediction problem now increases, since it requires learning a model that takes into account the semantics of the ontology in the output space, which is an additional requirement when compared to the current techniques that deal with general graphs or trees.

## 4. FORMALIZATION OF WEB BROWSING BEHAVIOR

When browsing the Web, users interact with Web resources via browsers interface (e.g. clicking links, submitting HTML forms, etc.). These interactions can be recorded as usage data in forms of Web server or client-side logs. We use the term *browsing event* to describe the basic component of user behavior in performing activities with the Web browser directly.

#### EXAMPLE 1. (Cross-site Browsing Logs)

ID Time	User Action
1 [17:04:14:35]	<a href="http://www.avis.com/car-rental/reservation/">http://www.avis.com/car-rental/reservation/</a>

start-reservation.ac?resForm.pickUpLocation=Lyon  
 1 [17:11:49:21] <http://www.google.de/search?q=Lyon+www2012>  
 1 [17:11:49:33] <http://dbpedia.org/page/Lyon>  
 1 [17:11:49:39] <http://data.semanticweb.org/conference/www/2011/demo/a-demo-search-engine-for-products>

In this running example of usage logs, the user starts a car rental reservation at Avis, next performs search at Google, and then visits sequentially the page of Lyon at DBpedia<sup>2</sup> and the page of a demo paper at Semantic Web Dog Food<sup>3</sup>. The last two sites are semantically-enabled, thus, offer a domain ontology and data publishing as Linked Open Data.

If each log entry is represented as a meaningful event from the application domain where it is issued, then user behavior becomes more comprehensible. The context of the event can be extended with additional information from the domain (as explained in Sec. 4.2).  $\square$

We aim at monitoring user behavior across multiple Web sites, logging clickthrough data upon agreement of Internet users. Each log entry is a tuple  $\mathcal{L} = \{UserID, URL, timestamp\}$  of a user ID, URL of the accessed Web resource, and real time when this happened. These usage data logs are initially stored in raw form, as produced upon each user interaction.

Our formalization approach, illustrated in Fig. 1, comprises a set of techniques, such as data preprocessing, human logs filtering, session construction, formal description of logs, and semantic enrichment. In the following sections, we present the definitions upon which this work is based. We then focus on the semantic enrichment approach, which extends our previous work [14].

The contributions of this work are: 1) a model to formally and semantically structure usage logs, 2) an approach for the automatic extraction of the semantics of usage logs from heterogeneous sources, in which the domain knowledge has a semi-structured or structured formal representation, and 3) a new approach for semi-supervised prediction with ontology-based output spaces

## 4.1 Formal Model for the Representation of Logs

**DEFINITION 1. (Event)** We define a browsing event as a tuple  $e = (l, \mathcal{T}, P, t)$ , where  $l$  is the full URL invoked,  $\mathcal{T}$  is a set of event types for which this event qualifies,  $P = \{p_1, \dots, p_l\}$  is a set of parameters and  $t$  is the occurrence time. For simplicity, we denote event time by  $e_i.t$  and set of event types by  $e_i.\mathcal{T}$ .

Each user browsing activity recorded in logs is physically represented by a URL, but conceptually it comprises an event that serves a particular function and relates to a specific content. We give meaning to each event issued as an HTTP request in the logs, by mapping its respective URL to domain concepts according to two dimensions: content and

function. An event resulting from the interaction of a user with a specific Web page serves a particular function (e.g. searching, browsing, login, etc.) related to some content (e.g. flight reservation, organization, hotel, etc.).

**DEFINITION 2. (Event Type)** An event can be mapped to several types denoted by the set  $\mathcal{T} = \{\mathcal{T}_c, \mathcal{T}_f\}$ , where  $\mathcal{T}_c$  is the type of content to which this event relates,  $\mathcal{T}_f$  is the type of function this event serves.

We have extended the definition of a browsing event with parameters, which can be extracted based on the information contained in the URL  $l$ . We consider three main conceptual elements in a link: URL base, variable names, and values. Based on the typical convention of URL formation, we syntactically split the link into two basic parts: URL base, which defines domain name, and the rest of the URL is used to extract input variables, which are modeled as event parameters.

**DEFINITION 3. (Parameter)** An event parameter  $p$ , which can be further classified as input or output parameter, is a pair  $p = (v_{name}, v_{value})$  consisting of variable name and value.

Events are grouped into sessions, which represent a period of sustained Web usage. The boundaries of a session are normally determined by temporal and behavioral factors (e.g. browsing intention). We follow previous research [5] in deploying an heuristic, which starts a new session after an idle period of 30 minutes between the browsing events.

**DEFINITION 4. (Session)** We denote a user session as a tuple  $S = \{s, T_s, T_e, U\}$ , where  $s = \langle e_1, e_2, \dots, e_n \rangle$  is an ordered sequence of browsing events performed from user  $U$ , such that  $e_i.t \leq e_{i+1}.t$  for all  $i$ , where  $i$  denotes the event order in the sequence. Furthermore,  $T_s$  is the starting time and  $T_e$  the ending time of the session, such that  $T_s \leq e_i.t \leq T_e$ .

Ordering of events in a session is used later as a feature for the supervised learning of event types, when they are not available in the domain ontology.

For the realization of these concepts, I have used a Web Browsing Activity Model (WAM), which I formalize as an ontology (Fig. 2). This is also presented in the paper Hoxha et al. [14]

**Classes and Properties.** Classes in WAM are divided in three groups: Core classes, External classes, and Type classes. External classes are basic concepts that I reuse from well-established ontologies. Each `wam:Event` is a subclass of the concept `event:Event` from the *Event ontology*<sup>4</sup>.

Each `wam:Session` has one `wam:StartEvent` and one `wam:EndEvent`, both of type `wam:Event`. Class `wam:User` is simply characterized by user IP address and ID, but the ontology allows

<sup>2</sup><http://dbpedia.org>

<sup>3</sup><http://data.semanticweb.org>

<sup>4</sup><http://purl.org/NET/c4dm/event.owl#>

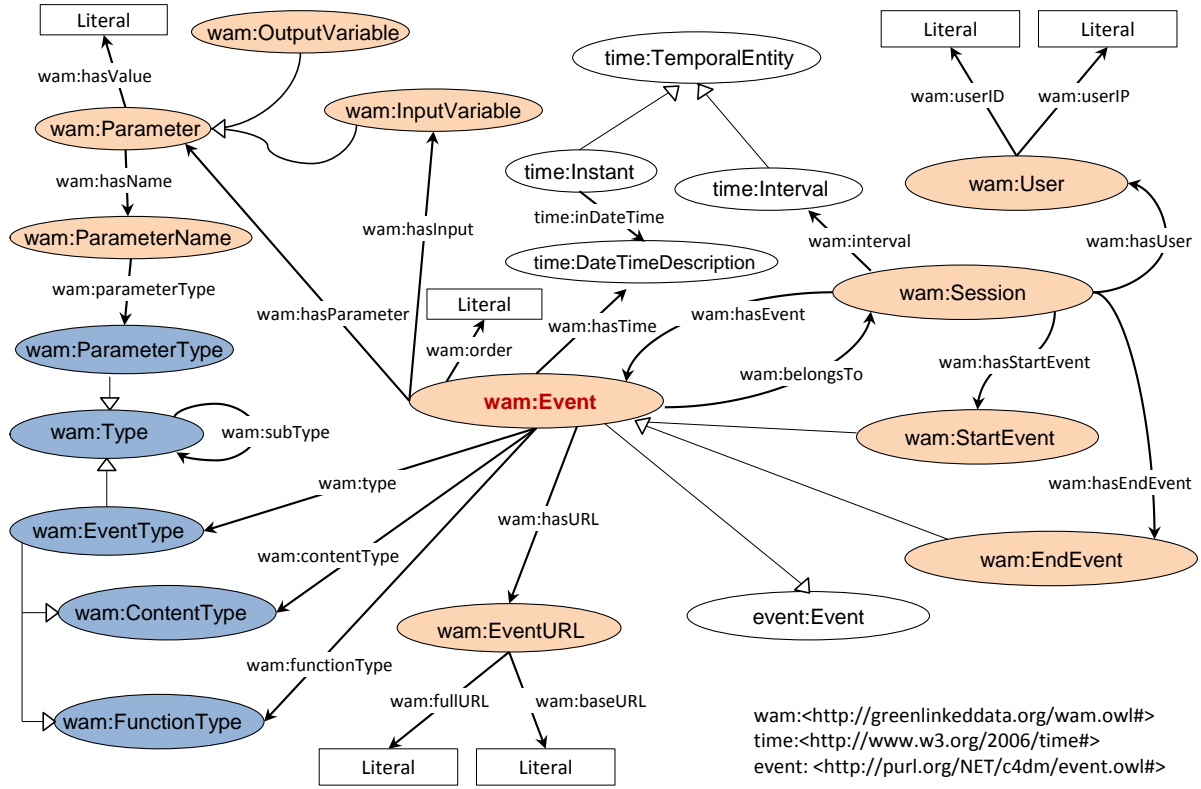


Figure 2: WAM Ontology

flexible future extendability with user profiles or other attributes (e.g IP-based geographical location). To annotate of event timestamps and session interval, I reuse basic concepts from *OWL Time ontology*,<sup>5</sup> which models knowledge about time such as temporal units, instants, etc.

The ontology is expressed in OWL-2-DL with underlying *SROIQ* logic [12].

## 4.2 Semantic Enrichment using Domain Knowledge

In order to obtain semantic information about the pages that the users have accessed, we extract domain-level structured objects as semantic resources contained in the pages.

The main focus of this semantic enrichment approach, formally described in algorithm Alg. 1, is to find the content types belonging to each event<sup>6</sup>. For each URL resource request  $l$  in the logs, we use the Content Negotiation mechanism<sup>7</sup> to retrieve the respective RDF representation, if it is available. Based on the RDF, we identify the Uniform Resource Identifier (URI) of the resource (Alg.1, line 4). We also retrieve the domain ontology  $O_d$  of the respective Web domain, as well as the class to which the given resource belongs (querying on *rdf:type*) (line 6). A resource may be a member of many classes in the  $O_d$ , therefore we consider all

of them as instances of content type (line 9).

**Algorithm 1** Automatic semantic enrichment of a browsing event: *findContentTypes(s)*

**Require:** Ordered sequence of events  $s = \langle e_1, e_2, \dots, e_n \rangle$   
**Ensure:** Update knowledge base with new ABox assertions  $\alpha_t$  related to content types

**for all**  $e_i \in s$  **do**  
    **Get the URL  $e_i.l$  of the event**  
3: Retrieve  $RDF_l = \text{retrieveRDF}(e_i.l)$  the RDF/XML representation  
    Get resource URI  $R_l = \text{identifyResourceURI}(RDF_l)$   
     $O_b = \text{getDomainOnt}(RDF_l)$ ;  
6: Find  $\mathcal{T}_c = \text{classMembership}(R_l)$   
    **for all**  $T \in \mathcal{T}_c$  **do**  
        Abox Assertions  
9:  $\alpha_t = \{ \text{ContentType}(T), \text{contentType}(e_i, T) \}$   
**Serialize assertions  $\alpha_t$  and update knowledge base**

An example of enriching an event with semantics is illustrated in Figure 3. This is the last log entry of our running example in Ex. 1. Initially, the raw log is just a syntactic representation of the URL request (in this case a demo paper). We retrieve the respective RDF representation and identify the resource with its URI. Querying (via SPARQL<sup>8</sup>) the domain ontology, here the *SWRC* publications ontology, we can enrich the event's semantics with additional knowledge. We find that this resource is a *Demo* of type *InProceedings*.

<sup>5</sup><http://www.w3.org/2006/time#>

<sup>6</sup>In this work, we focus on the contentType class only

<sup>7</sup><http://www.w3.org/Protocols/rfc2616/rfc2616-sec12.html>

<sup>8</sup><http://www.w3.org/TR/rdf-sparql-query/>

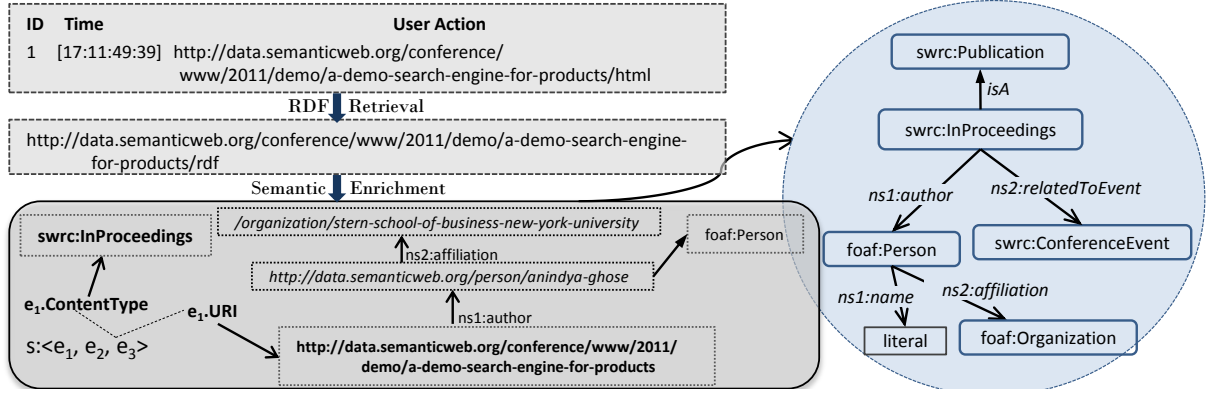


Figure 3: Semantic Enrichment of Usage Logs

We can further extend the context with information like the *conference* WWW2011 it belongs, the conference *location*, the *author* of the paper, the author's *affiliation*, etc.).

### 4.3 Leveraging Usage Logs with Structured Markup Data

The previous enrichment approach addressed Web sites that provide semantic annotations in pure RDF format, which even though increasing in popularity are still limited in number. For this reason, we also address another category of sites, which provide metadata as structured markup data, embedded in the HTML content. Annotation of HTML elements with structured markup data is a technique increasingly used nowadays by Web site providers [22] that enables search engines, web crawlers, and browsers to extract, automatically-process and better understand the content of pages. Structured markup data can be provided in different formats (such as RDFa, microdata, microformats), using particular supporting vocabularies (also called schemas), e.g. Open Graph Protocol,<sup>9</sup> schema.org,<sup>10</sup> DCMI Terms,<sup>11</sup> etc.

We devise an approach for the semantic enrichment of user logs with structured markup data, which we extract from the pages that the users have visited. The Web page under each URL from user logs is identified as a *Web resource* by a unique URI.<sup>12</sup> This resource is further enhanced with other objects, which annotate it with descriptive metadata based on shared schemas.

The semantic enrichment approach starts with the extraction of user browsing logs (tracked by a client toolbar) and segmenting logs in sessions, such that one session contains all requests of one user within a day. Therefore, there might also be different sessions that belong to the same user. We then filter those sessions that include pages belonging to several sites of interest. For our experiments, we have chosen a set of sites<sup>13</sup> from the events (concerts, conferences, etc.) advertisement domain.

The next steps consist in identifying the set of unique pages in the filtered user logs, then deploying metadata extraction and metadata analysis techniques. We map each page to a semantic Web resource, which we define (Def. 5) as the atomic unit of the modelling approach.

**DEFINITION 5. (Semantic Web Resource)** We define as a semantic Web Resource an information resource from the document Web that is identified by a dereferencable HTTP URI (Uniform Resource Identifier), and has an RDF/XML representation, which contains associated description of its attributes and relations to other Web resources.

It is important to highlight that different Web sites use different schemas to annotate their HTML elements. We semi-automatically align the concepts and relations among the schemas based on their respective semantics, in order to enable matching resources across different sites. The end result is a reference ontology  $\mathcal{O}$  consisting of constructs that define the concepts and their semantic relations used for the semantic annotation of resources across all sites. The resources are classified into *classes/concepts* of the ontology and are connected between each-other via semantic *relations* (e.g. *hasperformer*, *hasvideo*). Resources are also annotated by attributes, which are represented as RDF predicates and respective literal values.

We perform resource duplicate detection and entity linking, which is important for identifying pages that belong to different Web sites, but still semantically represent the same resource (e.g. same *performer*, *venue*, etc.).

We detect if two resources found under different links are duplicates, by aligning them based on their attribute predicates. We have manually identified a set of rules, which map predicates that belong to different schemas, but have the same semantics (e.g. <http://purl.org/dc/terms/title> and <http://opengraphprotocol.org/schema/title>).

We automatically group resources based on their type (ContentType), then compare the resources of the same type based on the values of the attributes, after aligning them by predicate name. The values of the attributes are compared using the Levenshtein distance. The Levenshtein distance

<sup>9</sup><http://ogp.me/>

<sup>10</sup><http://schema.org/>

<sup>11</sup><http://dublincore.org/documents/dcmi-terms/>

<sup>12</sup>Uniform Resource Identifier

<sup>13</sup>eventful.com, eventbrite.com and upcoming.yahoo.com



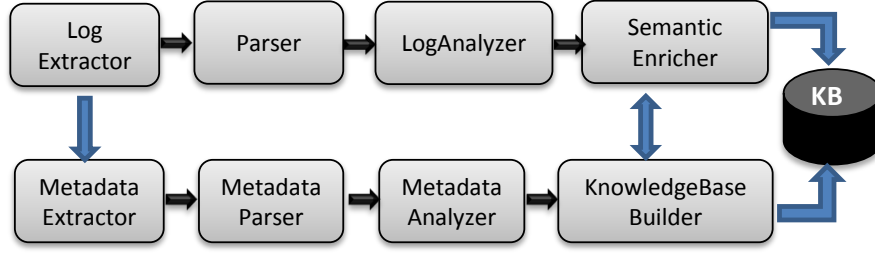


Figure 4: Leveraging Usage Logs with Structured Markup Data

between two strings is defined as the minimum number of edits needed to transform one string into the other, with the allowable edit operations being insertion, deletion, or substitution of a single character. In our case, we use a pre-define threshold value (0.8) of the edit distance between the attribute values. For example, two resources of type *Event* will be classified as duplicates if they have the same value of the attribute *venue* and *time*, whereas the values of their *title* attributes have an edit distance greater than the defined threshold.

## 5. LEARNING EVENT CONTENT TYPES

### 5.1 Problem Definition

Finding the **contentType** class of a browsing event can be formulated as a classification problem, borrowing from the field of machine learning. After the deployment of the formalization and automatic semantic enrichment approach of Sec. 4, we generate a session  $s_k = \langle e_1, e_2, \dots, e_n \rangle$  (Def. 4) as a sequence of semantically-annotated browsing events (as in Def. 1). For most of the events in the sequence, we are able to automatically find and assign a **contentType** class from the domain ontology. But, there are also two events in  $s_k$ , where no **contentType** class could be derived. Hence, we follow a second step of semantic enrichment that comprises a supervised technique for learning the class, based on the observed examples (i.e. already formalized events in the overall sessions). We need to assign a particular event from the logs to a predefined class, being in our case the **contentType** of this event. Hence, we approach our problem as a classification task, which learns a function  $f: E \rightarrow \mathcal{C}$  that maps an event  $e_i \in E$  s.t.  $e_i = (l_i, T_i, P_i, t_i)$  (as in Def. 1) to an output class  $c_i \in \mathcal{C}$ . In our case,  $\mathcal{C}$  is a set of classes belonging to an ontology  $\mathcal{O}$ .

### 5.2 Classification with Structural Support Vector Machines

In our approach, we use the generalized formulation of multi-class SVM learning [30]. We are interested on the problem of learning a function  $f: \mathcal{X} \rightarrow \mathcal{Y}$ , which maps input instances  $\mathbf{x} \in \mathcal{X}$ , which in our setting consist of the events in the logs, to discrete outputs  $\mathbf{y} \in \mathcal{Y}$  that consist of arbitrarily numbered labels representing **contentType** classes in our events ontology.

Lets consider the case of finding a function  $f$  that maps

each event  $\mathbf{x}_i$  from usage logs to one of the classes in  $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ . The problem addressed is to learn a discriminant function  $F: \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  over input/output pairs, so that for a given input  $\mathbf{x}$ , we can derive a prediction by maximizing  $F$  over the response variables:

$$F(\mathbf{x}; \mathbf{w}) = \underset{\mathbf{y} \in \mathcal{Y}}{\operatorname{argmax}} F(\mathbf{x}, \mathbf{y}; \mathbf{w}) \quad (1)$$

In our case, we deal with a multi-class classification problem [8], where  $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$  is the input set of events of sessions from the usage logs,  $\mathcal{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$  is the set of output classes from ontology  $\mathcal{O}$ , and  $\mathbf{w} = (w_1, \dots, w_N)$  is a stack of vectors with  $w_n$  being a weight vector for the class  $\mathbf{y}_n$ . We use the following formulations of the linear discriminant functions  $F$ :

$$F(\mathbf{x}, \mathbf{y}_n; \mathbf{w}) = \langle \mathbf{w}_n, \Phi(\mathbf{x}) \rangle \quad (2)$$

where  $\Phi(\mathbf{x}) \in \mathbb{R}$  is the vector of numeric features extracted from  $\mathbf{x}$ . SVM, then, solves the following optimization problem:

$$\min_{\mathbf{x}, \xi} \frac{1}{2} \sum_{i=1}^N \|\mathbf{w}_i\|^2 + \frac{C}{K} \sum_{i=1}^K \xi_i \quad (3a)$$

$$\forall i, \forall \mathbf{y} \in \mathcal{Y} \setminus \mathbf{y}_i : \langle \mathbf{w}, \Phi(x_i) \rangle \geq 100\Delta(\mathbf{y}_n, \hat{\mathbf{y}}_n) - \xi_i \quad (3b)$$

with regularization parameter  $C$  and slack variables  $\xi_i$  for margin violations (for details see [15]). The learning algorithm optimizes the error rate during training, minimizing prediction loss defined by a function  $\Delta: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}^D$ , where  $\Delta(\mathbf{y}_n, \hat{\mathbf{y}}_n)$  is the loss of predicting  $\mathbf{y}_n$  when the correct output is  $\hat{\mathbf{y}}_n$ .

#### Reasons for choosing structural SVM

There are several reasons for choosing structural SVM as our classification approach. Firstly, SVMs in general are shown to perform better in building complex and accurate models [15], particularly in settings similar to ours such as Web page categorization or purely URL-based page classification [17, 3]. Secondly, SVMs deal very well with sparse and highly dimensional data, as is the case of the huge and heterogeneous amounts of cross-site usage logs, which lead to feature vectors that are large and highly sparse. At last, structural SVMs enable learning for complex and interdependent objects of the output space, leading us towards

an extension of our approach in learning a formal, structured ontology with class relationships for the classification of events (i.e. requested resources) in the usage logs.

For our classification task, we follow a procedure comprising the following steps:

**Preparation of Training and Testing Datasets.** After the logs have been semantically formalized using our formalization approach, we select a portion of the data for the classification problem. Initially, since the formalized logs are represented as RDF triples and stored in a repository, using SPARQL queries we extract two sets of data for training and testing, each of them containing a huge vector of session id, URL of event and order of event belonging to that session. We then prepare training and test datasets, respectively. Since supervised learning needs labeled data, a part of those generated from the mapping to the domain ontology, which serve as ground truth values. The labels<sup>14</sup> that are not found in the ontology are annotated manually.

**Feature Selection.** We select different categories of features for the classification of event types. We first experiment with whole tokens (no stemming is performed) of URLs, and with the letter n-grams of the tokens [17]. We also test *sequential features*, such as sequences of pairs of tokens in the URL, referred as the *precedence bigrams*. We further propose a new feature based not only on the URL of the event, but on the sequential information related to the session in which the event belong (*sequential neighbors*). In this case, the tokens of the neighboring events are also included as features.

**Feature Vector Representation.** As explained earlier (Sec. 5.1), SVMs require that each instance in the input space is represented as a vector of real numbers. Hence, we convert our inputs into vectors of numeric values. In order to construct such feature vectors, we follow a series of preprocessing steps aligned with our definition of the features. Preprocessing includes tokenization, n-gram generation, and precedence bigram formation. Tokens or ngrams derived from the URL of the event serve as binary features.

**Model Selection.** We experiment with the linear kernel of structural SVM, motivated by the following reasons: high dimensionality of the feature vectors, huge number of features, and high number of classes/labels. We experiment with different values of the regularization parameter  $C$ .

We have conducted experiments with datasets of real-world usage logs. In section 6 we provide details on the characteristics of the datasets used for training and testing, as well as report on the evaluation results of these experiments.

## 6. EVALUATION

### 6.1 Formalization Approach

We provide a Java SE implementation of the introduced formalization approach, deploying the steps of processing usage logs, cleaning, and formalization with WAM ontology (whose consistency is checked with Pellet 1.5.2 reasoner).

<sup>14</sup>Terms *label* and *class*, as well as *instance* and *event* are used interchangeably.

We have further implemented the step of automatic semantic enrichment of events, for which we read and query using Jena Framework<sup>15</sup>.

In order to show the feasibility of this approach, we performed experiments by semantically formalizing logs from the USEWOD datasets [4] featured in Table 1. The formalized sessions and events are serialized in RDF, and then imported via OpenRDF Sesame Core 2.6.0 API<sup>16</sup> into a repository of a Sesame Framework<sup>17</sup> that is made available online<sup>18</sup>. Overall, we processed nearly one month of usage logs from large Web sites (such as DBPedia), proving that the approach is scalable and able to retrieve the content types classes of more than 80% of events.

Regarding the practicality of the proposed approach, which requires the existence of semantically-enabled websites or sites that include RDF annotations, we note that the percentage of such sites in the Web is now continuously and quickly increasing (for details see [22]).

### 6.2 Supervised Learning Approach

**Experimental Setup.** For our supervised learning experiments, we used two datasets  $D_1$  and  $D_2$  of different sizes extracted from the repository of events generated in the first step of our formalization. These are the events belonging to the two weeks of the SWDF part of Table 1. For both datasets we prepared training and testing sets. The test sets contain events for which the content type was not automatically found. We report on the characteristics of these datasets in Table 2.

For dataset  $D_1$ , we chose usage logs of two random consecutive days, extracting the events of one day (3. July) for the training and events of another day (2. July) for testing. Whereas for  $D_2$ , we chose a larger set comprising the logs of all the days from both weeks.

We use the implementation *structSVM*<sup>19</sup> of structural SVMs with the multi-class formulation. After experimenting with different values of the regularization parameter  $C$ , we chose the value 5000 to be the best one. For training, we follow a three-fold cross-validation approach.

**Evaluation Measures.** To evaluate the performance of our classification approach, we use the F-measure metric, which is the harmonic mean of precision ( $\pi$ ) and recall( $\rho$ ), defined as follows:

$$\pi_i = \frac{TP_i}{TP_i + FP_i}, \rho_i = \frac{TP_i}{TP_i + FN_i} \quad (4)$$

$$F_i = \frac{2\pi_i\rho_i}{\pi_i + \rho_i}, macroF_1 = \frac{\sum_{i=1}^N F_i}{N} \quad (5)$$

where  $TP_i$  (True Positives) is the number of instances assigned correctly to class  $i$ ;  $FP_i$  (False Positives) is the num-

<sup>15</sup><http://incubator.apache.org/jena/>

<sup>16</sup><http://www.openrdf.org/doc/sesame2/api/>

<sup>17</sup><http://www.openrdf.org/>

<sup>18</sup><http://46.4.66.131:8080/openrdf-workbench/repositories/wam/query>

<sup>19</sup>[http://svmlight.joachims.org/svm\\_struct.html](http://svmlight.joachims.org/svm_struct.html)



**Table 1: Results of the Formalization Approach**

	SWDF	DBPedia 3-3
Monitoring Period	01/07/09-13/07/09	01/07/09-13/07/09
Nr. Sessions	2831	31893
Average nr. sessions/day	235.92	2899
Nr. Triples	277788	> 3million
Nr. Events	10437	426k
Mode nr. events/session	4	10
% events with content-Type classes	83%	81%

**Table 2: Description of training and testing datasets**

	Dataset $D_1$		Dataset $D_2$	
	Training set	Test set	Training set	Test set
Nr. events/set	974	1152	4676	4957
Nr. events	2126		9632	
Nr. classes	66		82	

ber of instances that do not belong to class  $i$ , but are assigned to class  $i$  incorrectly; and  $FN_i$  (False Negatives) is the number of instances not assigned to class  $i$ , but which actually belong to this class.

The F-measure values are in the interval (0,1), and larger values correspond to higher classification quality. To compute the overall F-measure score of our multi-class classification problem, we use *macro-averaging* (Equation 5) as a binary evaluation measure across the overall  $N$  classes.

**Experimental Results.** In our experiments, we use the token feature as the baseline. All the other experiments additionally use each of the other features. We report on the zero/one-error (%) and macro-F1 measures of our results.

As can be observed in Table 3, in particular the ngram ( $N$ ) and sequential neighbor features ( $S$ ) play an important role in increasing the classification accuracy. For  $D_1$  we see that the combination of features  $N$  and  $S$  yields the most optimal results, since the error is the smallest, while still keeping a high value of the macro- $F_1$  measure (which is a harmony of precision and recall averaged across all classes)<sup>20</sup>. For  $D_2$  we note that the precedence bigram feature gives the best classification in terms of the 0/1 error rate. Still, as in  $D_1$ , the impact of the sequential neighbor feature yields the best combination of both the error and the overall averaged  $F_1$  score. This proves our expectation that users sequentially browse related resources, which can help us derive missing semantics.

## 7. APPLICATION: CROSS-DOMAIN RECOMMENDATIONS

<sup>20</sup>The experiments on  $D_2$ , whose results are reported as  $N/A$ , were not supported by our machine because of the high dimensionality of feature vectors.

We introduce a framework for computing top- $N$  recommendations to users while they are browsing different Web sites. The goal is to predict a set of Web resources that a user will be interested in visiting next. The objective of our approach is to recommend accurate, yet diverse resources that span across various domains. A cross-domain recommendation approach aims at predicting resources that are relevant to the user and not necessarily part of the same domain which she has visited or provided ratings for.

We introduce a two-step approach to generate recommendations. In the first step, we apply a machine learning approach to compute a set  $K$  of resources that are predicted as the most relevant for a user to visit next, given that she is at the moment at page  $i$  (referred as *query resource*). In the second step, to generate a final set of  $N$  recommendations that come from multiple and diverse domains, we apply a diversification approach on the set  $K$ , which also ensures to not compromise the accuracy of the pages predicted as highly relevant to the user.

From the Web navigation perspective, users often have browsing preferences that span across different pages and Web sites. Therefore, recommending pages belonging to diverse domains in terms of the type/category of the pages and the site where they are located, exposes the visitors to novel and possibly unexplored resources. Consider, for example, the scenario in Figure 5: suppose a Web visitor is viewing a page (referred as Web resource) on the site *eventful.com* about a *CMA Songwriter Series* (a country music event). A diversity-enhanced recommender system could recommend this visitor to view next not only musical events similar to this one, but also a page on the venue *Marathon Music Works* where it is organized, a page on the performer of this event, e.g. the country singer *Brantley Gilbert*, possibly in another Web site (*upcoming.yahoo.com*), a video of this performer in Youtube, and also similar performers of the same category (country singers), e.g. *Keith Urban*. The rationale behind our work is that semantic information of the Web resources, and especially knowledge of the relational structure inherent in their content, can play a significant role in providing cross-domain recommendations.

### 7.1 Relevance Model

The first step of our recommendation approach consists in finding which resources are relevant for a user to visit next, given that she is currently accessing a particular Web resource. We initially give our definition on resource *pair relevance* and *set relevance*. Later, we propose an approach to

**Table 3: Macro-F1 measure of the experimental results (Regularization parameter  $C = 5000$ )**

Feature Category	Dataset $D_1$			Dataset $D_2$		
	Nr. Features	zero/one-error (in %)	Macro-F1	Nr. Features	zero/one-error (in %)	Macro-F1
Token (T)	1357	14.40	0.79	4341	13.04	<b>0.75</b>
Trigram (N)	4673	14.41	<b>0.84</b>	11205	12.99	0.69
Precedence Bigram (P)	3385	14.75	0.82	11060	<b>11.80</b>	0.58
Sequential Neighbors (S)	4071	13.54	0.67	13023	<b>12.02</b>	<b>0.63</b>
S+P	6099	14.06	0.73	15647	N/A	N/A
N+S	7387	<b>13.45</b>	<b>0.74</b>	19887	N/A	N/A

estimate the values of resource pair relevance, and use these accordingly for generating a set  $K$  of recommendations.

**Pair Relevance.** Given two Web resources  $i$  and  $j$  as in Def. 5, let pair relevance  $P(rel|r_i, r_j)$  denote the quality value that captures the relevance of these resources to each other. We give a probabilistic interpretation to the quality values: they approximate the *likelihood* of resource  $j$  satisfying the user intent given the query resource  $r_i$ . In our case,  $P(rel|r_i, r_j)$  is determined by a scoring function based on user access patterns and content of resources  $i$  and  $j$ . Pair relevance can be seen as an item-item similarity measure, thus the order of resources in the pair is not important.

**Set Relevance.** Based on a probabilistic interpretation, we define the relevance of a set of recommendations as:

$$Rel(K|i) = 1 - \prod_{j \in K} (1 - P(rel|i, j)) \quad (1)$$

based on the following *independence* assumption: given a query resource  $i$ , the conditional probabilities of two other resources satisfying the user are independent. The probability that the user will find none of two resources resources  $j$  and  $l$  relevant equals  $(1 - P(rel|i, j))(1 - P(rel|i, l))$ , where the value  $(1 - P(rel|i, j))$  is the probability that  $j$  fails to satisfy. The probability that the set  $K$  will all fail to satisfy equals its product, by the independence assumption. One minus that product equals the probability that some resource in the set will satisfy the user.

The use of probabilities in our approach is motivated by the formulation of the set relevance and the joint objective function, defined in Sec. 7.4, for the maximization of diversity in the recommendation set.

## 7.2 Relevance Learning and Prediction

We formulate the problem of estimating resource pair relevance  $P(rel|r_i, r_j)$  as a binary classification task.

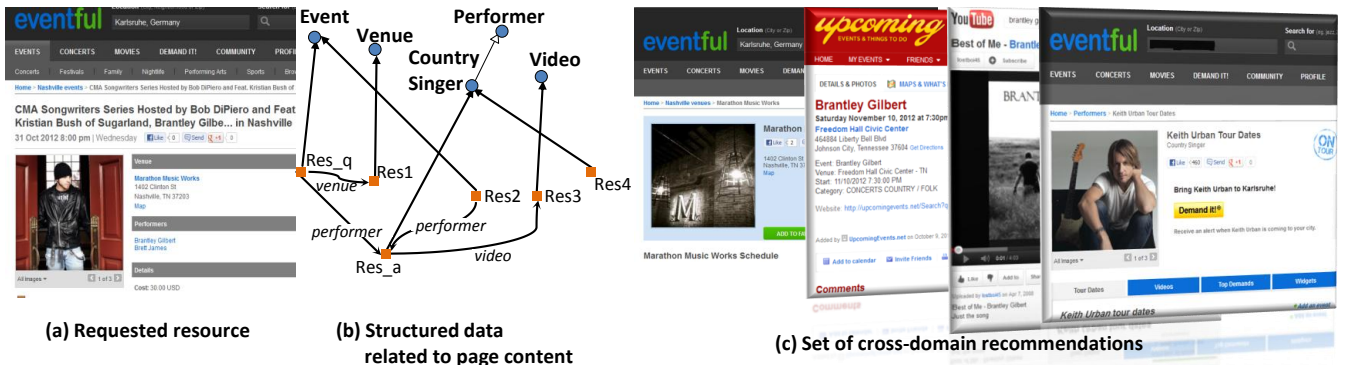
**Learning Pair Relevance.** The task is to learn a decision function  $f : \mathcal{R}^d \rightarrow Y$  based on an i.i.d training sample  $D_{train} = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)\}$ , where each training example consists of a feature vector  $\mathbf{x} \in \mathcal{R}^d$  and an output label  $y \in \{-1, 1\}$ . The learned function  $f$  is then used to predict the output label  $sign(f(\mathbf{x}_k))$  for test example  $\mathbf{x}_k$ .

Our original input data consist of a set  $X$  of resource pairs. Each pair  $x = \langle r_i, r_j \rangle \in X$  is initially mapped to a  $d$ -dimensional feature vector  $\mathbf{x}$  via a function  $\psi : X \rightarrow \mathcal{R}^d$ . The output labels in  $Y = \{-1, 1\}$  denote in our case the two classes: *non-relevant* and *relevant* resources in the pair.

**Probability Estimates.** For our relevance predictions, we are not just interested on hard decisions (labels), but rather the probability  $P(rel|r_i, r_j)$  (Eq. 1). We formulate it as an estimate of the confidence in the correctness of the predicted label. It is defined as the class conditional posterior probability  $P(y|\mathbf{x}) = P(y|\psi(x))$ , i.e. the probability with which the feature vector  $\mathbf{x}$  of pair  $x = \langle r_i, r_j \rangle$  belongs to class  $y$ . We deploy Support Vector Machines (SVMs) as probabilistic models by calibrating the scores into an accurate class conditional posterior probability with the sigmoid function [26]:

$$P(rel|r_i, r_j) = P(y = 1 | \psi(x = \langle r_i, r_j \rangle)) = \frac{1}{1 + \exp(Af(x) + B)} \quad (2)$$

fitted to the decision values of  $f$ , with parameters  $A$  and  $B$  estimated by minimizing the negative log likelihood of training data (using their labels and decision values) [18].



**Figure 5: Diverse Top-N Recommendations of Web Resources**

**Predicting Relevant Resources.** The approach allows us to learn a model, which we use to predict set  $K$  of resources that are the most relevant to a query resource  $r_i$ . We first derive pairs of resource  $r_i$  with other resources in the corpus, then apply the model learned with the SVMs to estimate the relevance value  $P(\text{rel}|r_i, r_j)$  for each pair. Afterwards, we select the top- $K$  resources with the highest relevance to  $r_i$ , ordering the pairs by the predicted pair relevance values. A crucial part of the prediction method is to define for the resource pairs the features (Sec. 7.3) that are effective in predicting an accurate relevance value.

### 7.3 Features

We engineer two groups of features: (1) usage-based features, which exploit the information contained in the user logs, and (2) content-based features that use the content of the resources. One feature especially captures the semantics of information in the ontology. Some of the features are<sup>21</sup>:

**SESSIONSIMILARITY:** a binary value stating if the two resources in the pair  $(r_i, r_j)$  appear together in at least one user session from the set  $S$ .

$$\text{sim}_{\text{session}}(r_i, r_j) = \begin{cases} 1, & \text{if } \text{Freqs}(r_i, r_j) > 0 \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

where  $\text{Freqs}(r_i, r_j)$  is the number of sessions in which resources  $r_i$  and  $r_j$  occur together.

**CONDITIONALSIMILARITY:** conditional probability of pair  $(r_i, r_j)$  occurring in the same session, given that resource  $r_i$  appears in that session:

$$\text{sim}_{\text{conditional}}(r_i, r_j) = \frac{\text{Freqs}(r_i, r_j)}{\text{Freqs}(r_i)} \quad (4)$$

**OBSERVEDRELEVANCEDEGREE:** Contrary to existing approaches that generally consider explicit user ratings as usage-based features, our challenging setting contains only implicit preference feedback inherent in user browsing logs. We introduce below a scheme to map implicit feedback into measurable values, in order to use them for identifying relevant resources as recommendations for users.

The measure of OBSERVEDRELEVANCEDEGREE (ORD) captures observations from usage patterns in the sessions of browsing logs. We model the correspondence between resource usage counts and user interest as a heuristic mapping between the access patterns and the probability of relevance. We adapt the Expected Reciprocal Rank metric [33, 7] for the setting of aggregated user sessions, introducing the scheme:

$$\text{ORD}(r_i, r_j) = \frac{2g(i,j)}{2g_{\text{max}}} \quad (5)$$

$$g(r_i, r_j) = n \cdot \mathcal{F}(\text{Freq}(r_i, r_j)) \quad (6)$$

$$\mathcal{F}(\text{Freq}(r_i, r_j)) = \frac{|\{r_k \in S' | \text{Freqs}(r_i, r_k) \leq \text{Freqs}(r_i, r_j)\}|}{|S'|} \quad (7)$$

<sup>21</sup>For more details see [13]

where  $S' \subseteq S$ . The value  $g(r_i, r_j)$  denotes the observed URL access frequencies in the overall user sessions. It is normalized to a common rating scale  $[0, n]$ , based on cumulative distribution function of  $\text{Freq}(r_i, r_j)$  over the set of other URLs accessed in the same session with  $r_i$  and  $r_j$ , but co-located with  $r_i$  less frequently. The maximum relevance value is  $g_{\text{max}}$ . In Equation (4), to compute the probability of relevance we do not subtract 1 as suggested in [7], in order to avoid overfitting to zero the probabilities of unobserved relevance [33].

**SYNTACTICSIMILARITY:** the term vector similarity measure between any two Web resources, computed as the cosine angle between two vectors modeled out of the *bag-of-words* (BOW) representation of the HTML page of each resource:

$$\text{sim}_{\text{syntactic}}(r_i, r_j) = \frac{\mathbf{V}(r_i) \cdot \mathbf{V}(r_j)}{|\mathbf{V}(r_i)| |\mathbf{V}(r_j)|} \quad (8)$$

$\mathbf{V}(i)$  is a real-valued vector composed of the weights of terms found in the HTML content of resource  $r_i$ . The weights are computed using the TF-IDF weighting scheme [20], product of Term Frequency (TF) and Inverse Document Frequency (IDF). The tf-idf weight is a statistical measure used to evaluate how important a word is to a document in a collection or corpus. In our case, we consider as document the HTML page of a resource. The content of the page is extracted, tokenization is performed, and the generated tokens are the words, whose weights are computed with the scheme. Resources are represented with the vector of tf-idf weights of tokens contained in their respective pages. As such, this feature entails only syntactic information.

**SEMANTICSIMILARITY:** a measure estimated via a set spreading approach [29] using the structural information related to the Web resources. Our spreading approach (Fig. 6) appends to a resource description terms that are related to the original terms based on an ontology. This is the ontology  $\mathcal{O}$  constructed in our semantic enrichment approach (Sec. 4.2).

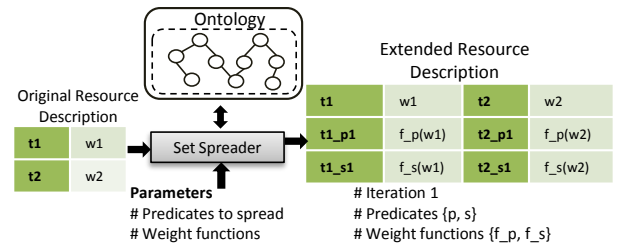


Figure 6: Set Spreading Approach

The process starts with an initial set  $RD_k = \{\langle t_k, w_k \rangle\}$  for each resource description, where  $t$  is the semantic type (e.g. *Performer*, *Event*, etc.) of this resource  $k$  denoted by a concept in the ontology  $\mathcal{O}$ , and weight  $w_k$  denotes the importance of the concept term in describing the resource. Each resource description  $RD_k$  is then iteratively extended via spreading, utilizing the concepts and relations (predicates) in  $\mathcal{O}$ . The set spreading of an  $RD_k$  results in the resource description  $RD'_k = \{\langle t_k, w_k \rangle, \langle t_{k-p}, f_{-p}(w_k) \rangle\}$ ,

which is extended by the term  $t_{k-p}$  related to  $t_k$  in  $\mathcal{O}$  by the predicate  $p$ . We pre-defined a set of predicates to find, at each iteration, *related* terms of the previous  $RDs$ . For example, the initial  $RD$  of a resource<sup>22</sup> annotated with type *Event* is the set  $\{\langle Event, 1.0 \rangle\}$ . Iteratively,  $RD$  is extended with the ontological terms related to the resource via predicates e.g. *hasVenue*, *title*, etc. We obtain next  $RD' = \{\langle Event, 1.0 \rangle, \langle B.B.KingBluesClub, 0.75 \rangle, \dots\}$ .

For the weights, we use a simple function  $f_p(w_k) = 0.75w_k$ , reducing the weights at each iteration<sup>23</sup>. The spreading process is terminated by predicates (relations) exhaustion.

The final similarity of two resources is then the mean cosine similarity of their descriptions  $RD_i$ . We compute the cosine similarity between two extended  $RDs$  at each iteration.

**SHARETYPE**: a binary value indicating if the pair  $(r_i, r_j)$  of resources have the same type (concept of the ontology  $\mathcal{O}$ ).

**SHARERELATION**: a binary value indicating if resources of pair  $(r_i, r_j)$  share a relation (in the ontology  $\mathcal{O}$ ) between them e.g. one resource is the event and the other resource is its venue, therefore sharing the relation *hasvenue*.

## 7.4 Diversity Model

While the issue of exploiting information from different domains to provide recommendations of items has been addressed from various perspectives, there has generally been an agreement among recent works [35, 1, 11] to address diversity in terms of topics or categories from a particular taxonomy. On this basis, we use the notion of item category in our definition of diversity. Furthermore, since we are dealing with resources of various Web sites, we also exploit the aspect of site diversity in our approach.

The item category is represented in our case by the *class type* with which a Web resource is semantically annotated. The definition of diversity (Def. 6) among Web resources covers two aspects: (1) semantic *type* to which the resources belong and (2) the *Web site* where they are located.

**DEFINITION 6. (DIVERSITY OF RESOURCES)** *Given a pair  $(i, j)$  of Web resources (as defined Def.5), they are considered diverse if any of the following two conditions occurs: (1) they have different semantic types, (2) they are located in different Web sites.*

For a measurable estimation, we define a distance function that considers both aspects of diversity in Def. 6. Specifically, given a set  $K$  of resources returned as *relevant* from the system for a given query resource  $i$ , to produce the final list of recommendations  $R$ , we measure the distance  $d(i, j)$  between any two resources  $j, l \in K$  by the function  $d : K \times K \rightarrow R$ , such that:

$$d(j, l; w_t) = w_t f_{type}(j, l) + (1 - w_t) f_{site}(j, l) \quad (9)$$

<sup>22</sup>Resource identifier URI <http://www.eventbrite.com/event/3504472973/>

<sup>23</sup>The rationale behind reducing weights is that the less directly related via terms in  $\mathcal{O}$  the resources are, the smaller their similarity degree should be at each iteration

$$s.t. \quad f_{type}(j, l) \in \{0, 1\}, f_{site}(j, l) \in \{0, 1\}, j, l \in K$$

where the distance function is symmetric  $d(j, l) = d(l, j)$  and computed as a weighted average of the values produced by functions  $f_{type}$  and  $f_{site}$ , which give 1 or 0 if *respectively* the *type* or the *site* of resources match, with weight  $w_t$  being the importance of diversity in *type*. Being application dependent, this weight can be configured to give more importance either to the diversity of types or to the diversity of Web sites, based on a particular desired scenario.

The overall diversity of a set of resources is modeled as the *average dissimilarity* of all pairs of resources contained in the set. Specifically, we use the averaged intra-list distance (ILD) metric [35, 33]:

$$Div(R) = \frac{1}{|R|(|R| - 1)} \sum_{j \in R} \sum_{l \in R, j \neq l} d(j, l) \quad (10)$$

where, in our case  $d(i, j)$  is the distance function of Eq. 9.

## 7.5 Diversity Enhancement

The initial set of recommendations is constructed with a relevance maximization method that follows a similarity-based approach (Sec. 7.1). As such, the rationale behind enhancing the diversity of the recommendation set is that, the resources generated for recommendation are very likely to be similar to each other. If the final recommendations comprise a diverse set of Web resources, it is more likely that a user finds in this set relevant items that fulfill her navigation intent.

However, an approach that displays to the user diverse, but non-relevant resources is not able to offer satisfactory results. Therefore, the goal of jointly offering a final set  $R$  recommendations with high diversity and a set  $R$  of high (similarity-based) relevance are opposite to each-other. To address this issue, we propose an approach to enhance diversity by maximizing an objective function, which captures the trade-offs between the relevance and diversity of the recommendations set. Our diversity maximization objective is the following:

Given a query resource  $r_i$ , a set  $K$  of resources qualified as relevant to  $r_i$ , an integer  $N$ , and a fixed control parameter  $\lambda \in [0, 1]$ , find the set of resources  $R \subseteq K$  with  $|R| = N$  that maximize the following combined objective function, representing a trade-off between diversity and relevance of the set.

$$f_t(K, N, \lambda) \triangleq (1 - \lambda) \alpha \mathbf{Rel}(R|i) + \lambda \beta \mathbf{Div}(R) \quad (11)$$

The objective in Equation (11) is the weighted arithmetic mean of set relevance and set diversity, where  $\alpha$  and  $\beta$  are normalization parameters that ensure these two measures are normalized to the same scale in the range  $[0, 1]$ . The parameter  $\lambda$  controls the degree of trade-off.

### 7.5.1 Diversity Maximization Algorithm

Our approach follows the set selection problem, finding the best set that maximizes the objective function and then ranking the set in order of *pair relevance* between each set element and the query resource.

Given the parameters  $K$ ,  $\lambda \in [0, 1]$ , and a given integer  $N \in \mathbb{Z}^+$  the objective is to select a set of resources  $R \subseteq K$ , such that the value of the function  $f_t$  (Eq. 11) is maximized, i.e. the objective is to find

$$R^* = \operatorname{argmax}_{R \subseteq K, |R|=N} f_t(K, N, \lambda) \quad (12)$$

s.t. all arguments besides  $R$  are fixed inputs of the function.

---

**Algorithm 2** Diversity maximization algorithm

---

**Require:** Set  $K$ , control parameter  $\lambda$ , integer  $N < |K|$

**Ensure:** Set  $R$  ( $|R| = N$ ) that maximizes  $f_t$

Initialize the set  $S = \emptyset$ ;

Find subsets of  $K$  with length  $N$ ;

3: for each subset, find  $f_t(K, N, \lambda)$   
output the set  $R$  with the highest value

---

We apply a brute force solution strategy, depicted in Algorithm 1, which creates the power subsets of length  $N$  from the resources in set  $K$ . For each subset, we calculate the value of the function  $f_t$  and return the subset with the maximal value as the final output of recommendations. Note that the set  $K$  is generated with the resources from the pairs estimated as the most relevant, based on our prediction approach. As such, the size of  $K$  need not be very big in order to target the highest relevant resources. While there are other methods that can be also used to solve this constrained maximization problem, we find this strategy effective in finding a local optimum for these values of  $|K|$  in our diversity enhancement setting.

## 8. CONCLUSIONS AND FUTURE WORK

In this paper, we address the issues of modeling user Web browsing behavior, tackling the problems related to information heterogeneity by using semantics. We present an approach for the formalization of such behavior across multiple sites based on a newly-introduced Web browsing Activity Model (WAM). A crucial part of the formalization is a two-staged semantic enrichment of logs, which maps them to events with comprehensible content types from the application domain. In order to find such semantic annotations of the logs, in the first stage we perform an automatic technique to retrieve the semantic types of logs from existing domain ontologies of Web sites. To annotate the remaining logs of those sites that do not provide a formal domain ontology, we deploy a supervised learning technique via a multi-class classification formulation. We explore for the first time the use of Support Vector Machines with structural and inter-dependent output spaces, as well as the exploration of new session-related sequential features for the semantic classification of usage logs.

The semantically-leveraged logs provide an added-value with respect to their syntactic representation in various ways: allow for more expressive formulation of queries to discover user navigation patterns; are useful input for techniques, such as semantic pattern mining, next-step navigation prediction or user clustering, which usually assume that these semantics of logs exists or are manually derived. A more beneficial aspect is the extension of these techniques to deal with cross-site browsing data and not only single Web sites.

We have implemented the overall formalization approach

with both stages of the semantic enrichment and performed experiments with real-world datasets of usage logs. We show that the extension with the supervised classification technique increases considerably the annotation accuracy. The introduced sequential features play an important role in ensuring a higher classification quality.

We plan to further investigate the techniques of learning semantic types (in particular function types) of usage logs, especially for semi-supervised techniques that reduce the effort of manually labeling training data. More interestingly, this work lays the foundations for a promising learning problem where the output space is not a set of classes, but a structured, formal ontology containing also the relations among concepts. We will elaborate on these aspects in our future work.

## 9. REFERENCES

- [1] R. Agrawal, S. Gollapudi, A. Halverson, and S. Ieong. Diversifying search results. In *Proceedings of the Second ACM International Conference on Web Search and Data Mining, WSDM '09*, pages 5–14, New York, NY, USA, 2009. ACM.
- [2] C. Baier and J.-P. Katoen. *Principles of Model Checking*. The MIT Press, 2008.
- [3] E. Baykan, M. Henzinger, L. Marian, and I. Weber. A comprehensive study of features and algorithms for url-based topic classification. *TWEB*, 5(3):15, 2011.
- [4] B. Berendt, L. Hollink, V. Hollink, M. Luczak-Rösch, K. H. Möller, and D. Vallet, editors. *USEWOD2012-2nd International Workshop on Usage Analysis and The Web of Data*, Lecture Notes in Computer Science. Springer, 2012.
- [5] R. E. Bucklin and C. Sismeiro. A Model of Web Site Browsing Behavior Estimated on Clickstream Data. *Journal of Marketing Research*, XL:249–267, Aug. 2003.
- [6] L. D. Catledge and J. E. Pitkow. Characterizing browsing strategies in the world-wide web. In *Computer Networks and ISDN Systems*, pages 1065–1073, 1995.
- [7] O. Chapelle, D. Metzler, Y. Zhang, and P. Grinspan. Expected reciprocal rank for graded relevance. In *Proceedings of the 18th ACM conference on Information and knowledge management, CIKM '09*, pages 621–630, New York, NY, USA, 2009. ACM.
- [8] K. Crammer, Y. Singer, N. Cristianini, J. Shawe-taylor, and B. Williamson. On the algorithmic implementation of multiclass kernel-based vector machines. *Journal of Machine Learning Research*, 2:2001, 2001.
- [9] M. d’Áquin, S. E. L., and E. Motta. Semantic technologies to support the user-centric analysis of activity data. In *Workshop on Social Data on the Web Workshop (SDoW) at ISWC*, 2011.
- [10] D. Downey, S. T. Dumais, and E. Horvitz. Models of searching and browsing: Languages, studies, and application. In *Proceedings of IJCAI*, pages 2740–2747, 2007.
- [11] I. Fernández-Tobías, I. Cantador, M. Kaminskis, and F. Ricci. A generic semantic-based framework for cross-domain recommendation. In *Proceedings of the*

- 2nd International Workshop on Information Heterogeneity and Fusion in Recommender Systems*, HetRec '11, pages 25–32, New York, NY, USA, 2011. ACM.
- [12] I. Horrocks, O. Kutz, and U. Sattler. The even more irresistible *SRIOQ*. In *Proceedings of the 10th International Conference on Principles of Knowledge Representation and Reasoning (KR2006)*, pages 57–67, June 2006.
  - [13] J. Hoxha. Semantic formalization of cross-site user browsing behavior. In *Research Technical Report, Archiv nr. 3025*. Institut AIFB, KIT, <http://www.aifb.kit.edu/web/Techreport3025/en>, 2012.
  - [14] J. Hoxha, M. Junghans, and S. Agarwal. Enabling semantic analysis of user browsing patterns in the web of data. In *USEWOD Workshop at the 21st International World Wide Web Conference (WWW2012)*, volume abs/1204.2713, 2012.
  - [15] T. Joachims, T. Hofmann, Y. Yue, and C.-N. J. Yu. Predicting structured objects with support vector machines. *Commun. ACM*, 52(11):97–104, 2009.
  - [16] E. J. Johnson, W. W. Moe, P. S. Fader, S. Bellman, and G. L. Lohse. On the depth and dynamics of online search behavior. *Manage. Sci.*, 50:299–308, March 2004.
  - [17] M.-Y. Kan and H. O. N. Thi. Fast webpage classification using url features. In *Proceedings of the 14th ACM international conference on Information and knowledge management*, CIKM '05, pages 325–326, New York, NY, USA, 2005. ACM.
  - [18] H.-T. Lin, C.-J. Lin, and R. C. Weng. A note on platt's probabilistic outputs for support vector machines. *Mach. Learn.*, 68(3):267–276, Oct. 2007.
  - [19] N. R. Mabroukeh and C. I. Ezeife. Using domain ontology for semantic web usage mining and next page prediction. In *CIKM*, pages 1677–1680, 2009.
  - [20] C. D. Manning, P. Raghavan, and H. Schütze. Introduction to information retrieval. *Information Retrieval*, 13(2):192–195, Apr. 2010.
  - [21] A. Mehdi, P. Valtchev, R. Missaoui, and C. Djeraba. A framework for mining meaningful usage patterns within a semantically enhanced web portal. In B. C. Desai, C. K.-S. Leung, and S. P. Mudur, editors, *C3S2E*, ACM International Conference Proceeding Series, pages 138–147, 2010.
  - [22] P. Mika and T. Potter. Metadata statistics for a large web corpus. In *Proceedings of the Linked Data Workshop (LDOW) at the International World Wide Web Conference*, 2012.
  - [23] A. L. Montgomery and C. Faloutsos. Identifying web browsing trends and patterns. *Computer*, 34:94–95, July 2001.
  - [24] D. Oberle, B. Berendt, A. Hotho, and J. Gonzalez. Conceptual user tracking. In E. M. Ruiz, J. Segovia, and P. S. Szczepaniak, editors, *AWIC*, volume 2663 of *Lecture Notes in Computer Science*, pages 155–164. Springer, 2003.
  - [25] Y. H. Park and P. S. Fader. Modeling browsing behavior at multiple websites. *Marketing Science*, pages 280–303, 2004.
  - [26] J. Platt. Probabilistic outputs for support vector machines and comparison to regularized likelihood methods. In A. J. Smola, P. Bartlett, B. Schoelkopf, and D. Schuurmans, editors, *Advances in Large Margin Classifiers*, pages 61–74. MIT Press, 2000.
  - [27] R. Stühmer, D. Anicic, S. Sen, J. Ma, K.-U. Schmidt, and N. Stojanovic. Lifting events in rdf from interactions with annotated web pages. In *Proceedings of the 8th International Semantic Web Conference*, ISWC '09, pages 893–908. Springer-Verlag, 2009.
  - [28] G. Stumme, A. Hotho, and B. Berendt. Usage mining for and on the semantic web. In *Next Generation Data Mining. Proc. NSF Workshop*, pages 77–86, Baltimore, 2002.
  - [29] R. Thiagarajan, G. Manjunath, and M. Stumptner. Computing semantic similarity using ontologies. In *Technical Report*. HP, 2008.
  - [30] I. Tsochantaridis, T. Hofmann, T. Joachims, and Y. Altun. Support vector machine learning for interdependent and structured output spaces. In *Proceedings of the 21. International Conference on Machine learning*, NY, USA, 2004. ACM.
  - [31] M. Tvarozek, M. Barla, and M. Bieliková. Personalized presentation in web-based information systems. In *Proceedings of the 33rd conference on Current Trends in Theory and Practice of Computer Science*, pages 796–807, Berlin, Heidelberg, 2007. Springer-Verlag.
  - [32] M. Vanzin, K. Becker, and D. D. A. Ruiz. Ontology-based filtering mechanisms for web usage patterns retrieval. In *EC-Web'05*, pages 267–277, 2005.
  - [33] S. Vargas and P. Castells. Rank and relevance in novelty and diversity metrics for recommender systems. In *Proceedings of the fifth ACM conference on Recommender systems*, RecSys '11, pages 109–116, New York, NY, USA, 2011. ACM.
  - [34] H. Yilmaz and P. Senkul. Using ontology and sequence information for extracting behavior patterns from web navigation logs. In *Data Mining Workshops (ICDMW), 2010 IEEE International Conference on*, pages 549–556, dec. 2010.
  - [35] C.-N. Ziegler, S. M. McNee, J. A. Konstan, and G. Lausen. Improving recommendation lists through topic diversification. In *Proceedings of the 14th international conference on World Wide Web*, WWW '05, pages 22–32, New York, NY, USA, 2005. ACM.